

# BAYESIAN DIVERGENCE-TIME ESTIMATION

**Tracy Heath**

Ecology, Evolution, & Organismal Biology  
Iowa State University



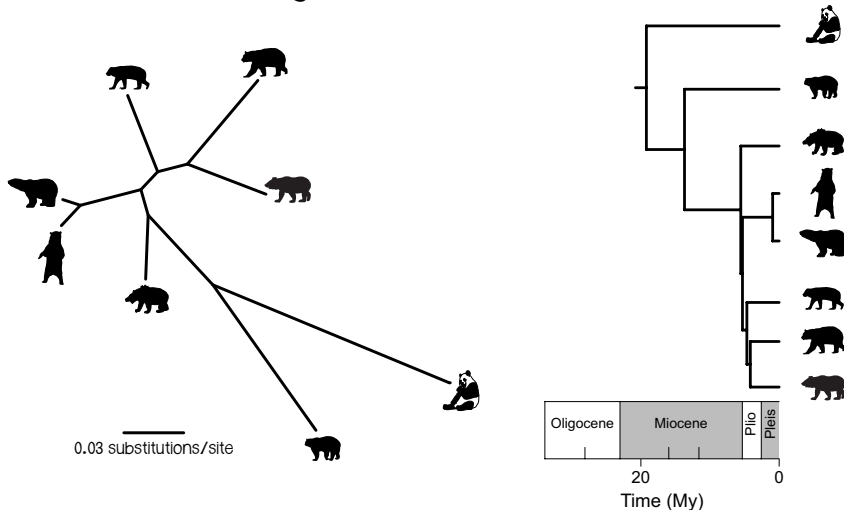
@tracy7

<http://phyloworks.org>

Workshop on Phylogenomics  
Český Krumlov, CZ  
January 28, 2019

# A TIME-SCALE FOR EVOLUTION

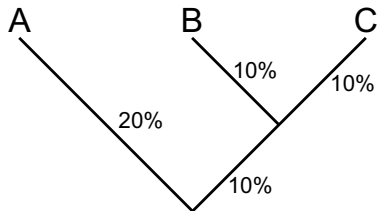
Phylogenies with branch lengths proportional to time provide more information about evolutionary history than unrooted trees with branch lengths in units of substitutions/site.



# THE GLOBAL MOLECULAR CLOCK

Assume that the rate of evolutionary change is constant over time

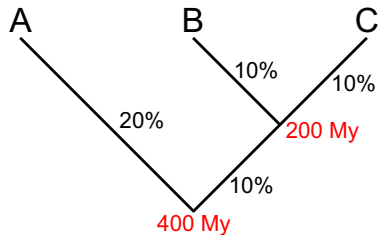
(branch lengths equal percent sequence divergence)



(Based on slides by Jeff Thorne; <http://statgen.ncsu.edu/thorne/compmolevo.html>)

# THE GLOBAL MOLECULAR CLOCK

We can date the tree if we know the rate of change is 1% divergence per 10 My

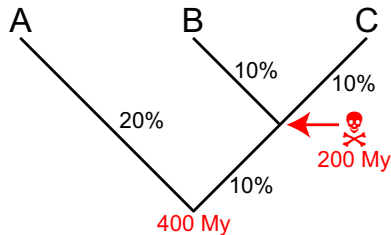


(Based on slides by Jeff Thorne; <http://statgen.ncsu.edu/thorne/compmolevo.html>)



# THE GLOBAL MOLECULAR CLOCK

If we found a fossil of the MRCA of **B** and **C**, we can use it to calculate the rate of change & date the root of the tree



(Based on slides by Jeff Thorne; <http://statgen.ncsu.edu/thorne/compmolevo.html>)

# REJECTING THE GLOBAL MOLECULAR CLOCK

Rates of evolution vary across lineages and over time

## **Mutation rate:**

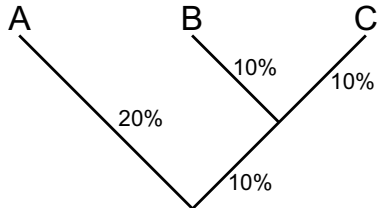
Variation in

- metabolic rate
- generation time
- DNA repair

## **Fixation rate:**

Variation in

- strength and targets of selection
- population sizes

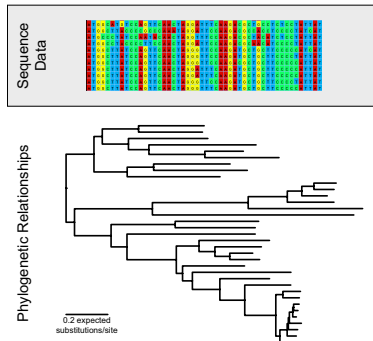


# UNCONSTRAINED ANALYSIS

Sequence data provide information about **branch lengths**

In units of **the expected # of substitutions per site**

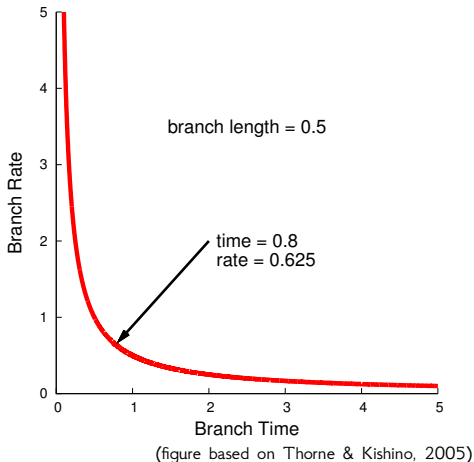
branch length = rate  $\times$  time



# ESTIMATING RATE & TIME

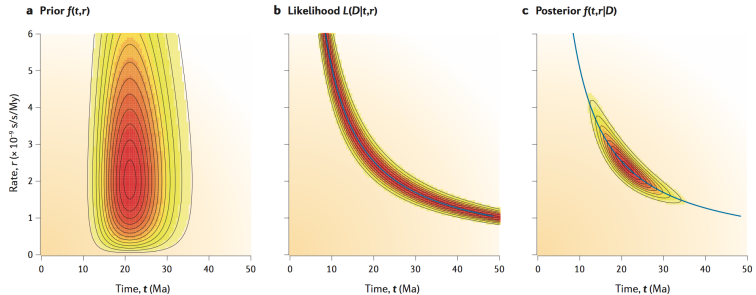
The sequence data provide information about branch length

for any possible rate, there's a time that fits the branch length perfectly



# ESTIMATING RATE & TIME

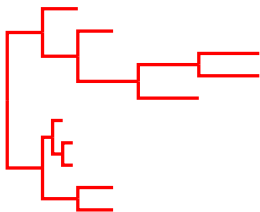
Methods for dating species divergences estimate the **substitution rate** and **time** separately



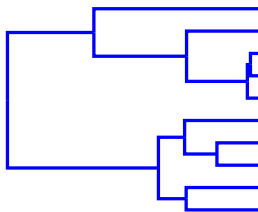
(dos Reis et al. *Nature Reviews Genetics*, 2016)

Tree-time priors for molecular phylogenies are only informative on a **relative** time scale

# BAYESIAN DIVERGENCE TIME ESTIMATION



length = rate



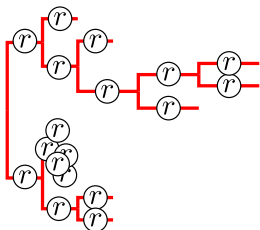
length = time

$$\mathcal{R} = (r_1, r_2, r_3, \dots, r_{2N-2})$$

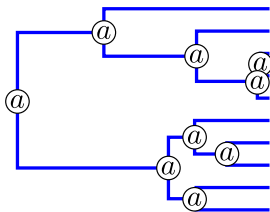
$$\mathcal{A} = (a_1, a_2, a_3, \dots, a_{N-1})$$

$$N = \text{number of tips}$$

# BAYESIAN DIVERGENCE TIME ESTIMATION



length = rate



length = time

$$\mathcal{R} = (r_1, r_2, r_3, \dots, r_{2N-2})$$

$$\mathcal{A} = (a_1, a_2, a_3, \dots, a_{N-1})$$

$$N = \text{number of tips}$$

# BAYESIAN DIVERGENCE TIME ESTIMATION

Posterior probability

$$f(\mathcal{R}, \mathcal{A}, \theta_{\mathcal{R}}, \theta_{\mathcal{A}}, \theta_s, \mathcal{T} \mid D)$$

$\mathcal{R}$	Vector of rates on branches
$\mathcal{A}$	Vector of internal node ages
$\theta_{\mathcal{R}}, \theta_{\mathcal{A}}, \theta_s$	Model parameters
$D$	Molecular or morphology data
$\mathcal{T}$	Tree topology



# BAYESIAN DIVERGENCE TIME ESTIMATION

$$f(\mathcal{R}, \mathcal{A}, \theta_{\mathcal{R}}, \theta_{\mathcal{A}}, \theta_s | D) =$$

$$\frac{f(D | \mathcal{R}, \mathcal{A}, \theta_s) f(\mathcal{R} | \theta_{\mathcal{R}}) f(\mathcal{A} | \theta_{\mathcal{A}}) f(\theta_s)}{f(D)}$$

$$f(D | \mathcal{R}, \mathcal{A}, \theta_s)$$

Likelihood

$$f(\mathcal{R} | \theta_{\mathcal{R}})$$

Prior on rates

$$f(\mathcal{A} | \theta_{\mathcal{A}})$$

Prior on node ages

$$f(\theta_s)$$

Prior on substitution parameters

$$f(D)$$

Marginal probability of the data

# MODELING RATE VARIATION

Some models describing lineage-specific substitution rate variation:

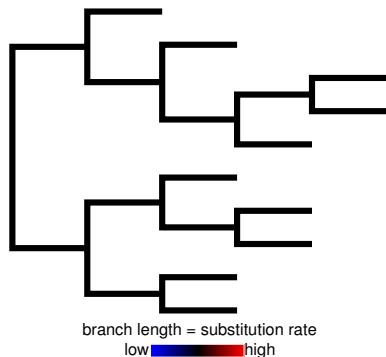
- **Global/strict clock** (Zuckermandl & Pauling, 1962)
- **Local clocks** (Hasegawa, Kishino & Yano 1989; Kishino & Hasegawa 1990; Yoder & Yang 2000; Yang & Yoder 2003, Drummond and Suchard 2010)
- **Punctuated rate change model** (Huelsenbeck, Larget and Swofford 2000)
- **Autocorrelated rates** (Thorne, Kishino & Painter 1998; Kishino, Thorne & Bruno 2001; Thorne & Kishino 2002; Lepage et al. 2007)
- **Mixture models on branch rates** (Heath, Holder, Huelsenbeck 2012)
- **Uncorrelated/independent rates models** (Drummond et al. 2006; Rannala & Yang 2007; Lepage et al. 2007)

# GLOBAL CLOCK

The substitution rate is constant over time

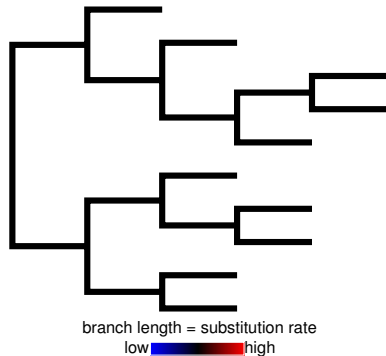
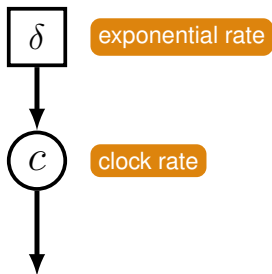
All lineages share the same rate

(Zuckerkandl & Pauling, 1962)



# GLOBAL CLOCK

$$c \sim \text{Exponential}(\delta)$$

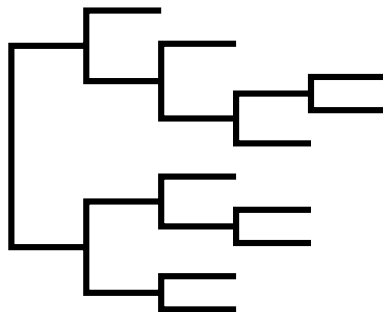
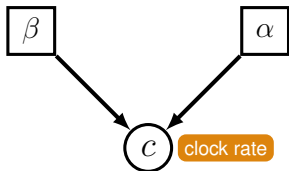



# GLOBAL CLOCK

$$c \sim \text{Gamma}(\alpha, \beta)$$

gamma scale

gamma shape



branch length = substitution rate  
low  high

# RELAXED-CLOCK MODELS

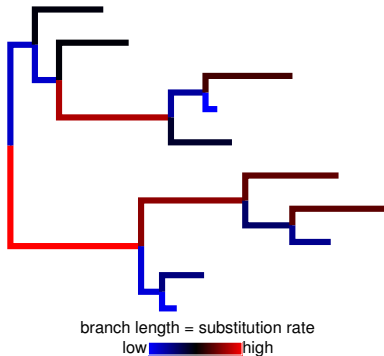
To accommodate variation in substitution rates  
'relaxed-clock' models estimate lineage-specific substitution rates

- **Local clocks**
- **Punctuated rate change model**
- **Autocorrelated rates**
- **Mixture models on branch rates**
- **Uncorrelated/independent rates models**

# INDEPENDENT/UNCORRELATED RATES

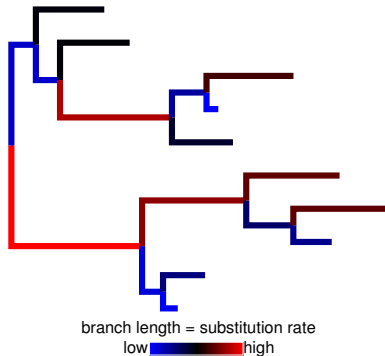
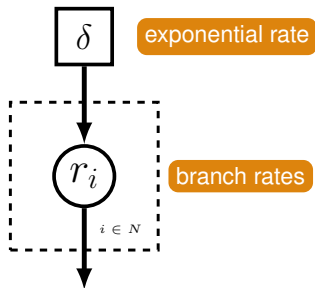
Lineage-specific rates are uncorrelated when the rate assigned to each branch is independently drawn from an underlying distribution

(Drummond et al. 2006; Rannala & Yang 2007; Lepage et al. 2007)



# INDEPENDENT/UNCORRELATED RATES

$$r_i \sim \text{Exponential}(\delta)$$



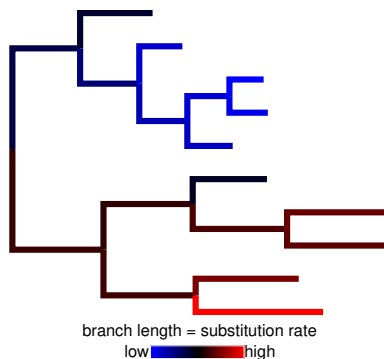


# AUTOCORRELATED RATES

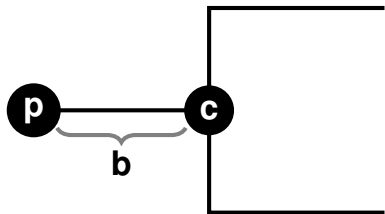
Substitution rates evolve gradually over time – closely related lineages have similar rates

The rate at a node is drawn from a distribution with a mean equal to the parent rate

(Thorne, Kishino, Painter, 1998;  
Kishino, Thorne, Bruno, 2001)



# AUTOCORRELATED RATES

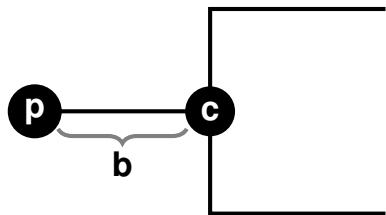


p = parent node

c = child node

b = branch

# AUTOCORRELATED RATES



p = parent node

c = child node

b = branch

$$r_c \sim \text{Lognormal}(\mu_c, \sigma_c)$$

$$\sigma_c := \nu t_b$$

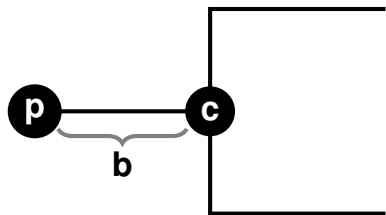
$$\mu_c := \ln(r_p) - \frac{\sigma_c^2}{2}$$

$$r_b := \frac{r_p + r_c}{2}$$

$\nu$  = variance parameter

$t_b$  = time duration of branch

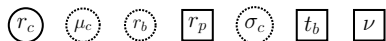
# AUTOCORRELATED RATES



p = parent node

c = child node

b = branch



$$r_c \sim \text{Lognormal}(\mu_c, \sigma_c)$$

$$\sigma_c := \nu t_b$$

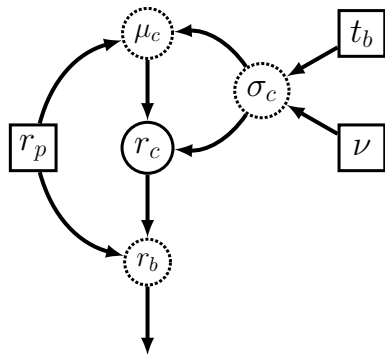
$$\mu_c := \ln(r_p) - \frac{\sigma_c^2}{2}$$

$$r_b := \frac{r_p + r_c}{2}$$

$\nu$  = variance parameter

$t_b$  = time duration of branch

# AUTOCORRELATED RATES



$$r_c \sim \text{Lognormal}(\mu_c, \sigma_c)$$

$$\sigma_c := \nu t_b$$

$$\mu_c := \ln(r_p) - \frac{\sigma_c^2}{2}$$

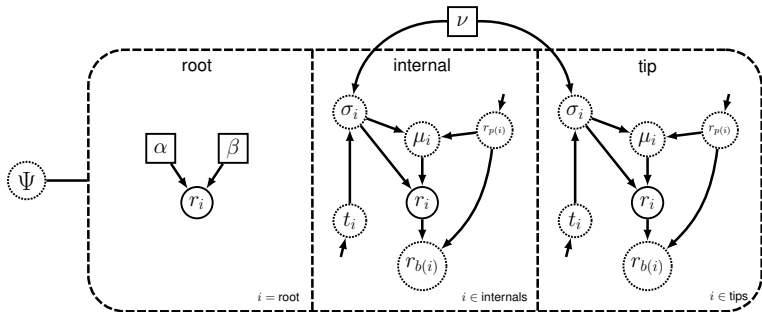
$$r_b := \frac{r_p + r_c}{2}$$

$\nu$  = variance parameter

$t_b$  = time duration of branch

# AUTOCORRELATED RATES

The rate associated with each node is a stochastic node, drawn from a distribution centered on its parent node



There is a gamma prior distribution on the rate at the root node

# MODELING RATE VARIATION

These are only a subset of the available models for branch-rate variation

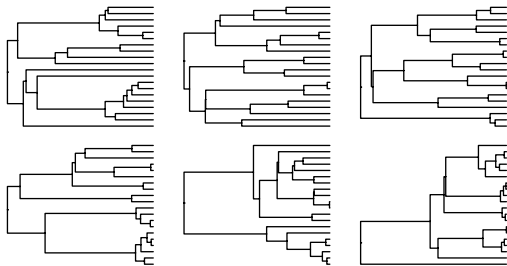
- **Global/strict clock**
- **Local clocks**
- **Punctuated rate change model**
- **Log-normally distributed autocorrelated rates**
- **Mixture models on branch rates**
- **Uncorrelated/independent rates models**

Considering model selection, uncertainty, & plausibility is **very** important for Bayesian divergence time analysis



# PRIORS ON THE TREE AND NODE AGES

Relaxed clock Bayesian analyses require a prior distribution on time trees



Different tree priors make different assumptions about the timing of divergence events and shape of the tree topology



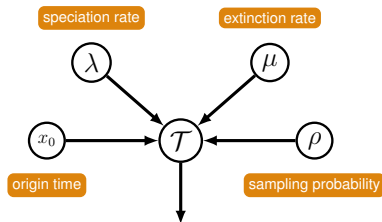
# STOCHASTIC BRANCHING PROCESSES

Tree priors based on stochastic models of lineage diversification

## Birth-death-sampling

**process:** at any point in time a lineage can speciate at rate  $\lambda$  or go extinct with a rate of  $\mu$

Conditions on a probability of sampling a tip,  $\rho$  and the origin time of the process,  $\varphi$

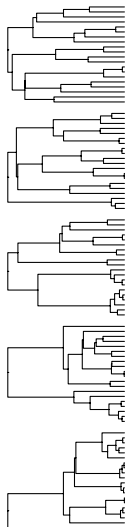


# STOCHASTIC BRANCHING PROCESSES

Different values of  $\lambda$  and  $\mu$  lead to different trees

Bayesian inference under these models can be very sensitive to the values of these parameters

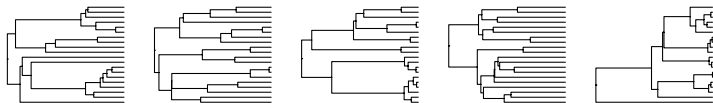
Using hyperpriors on  $\lambda$  and  $\mu$  (or  $d$  and  $r$ ) accounts for uncertainty in these hyperparameters



# PRIORS ON THE TREE AND NODE AGES

Sequence data are only informative on *relative* rates & times

Most tree priors cannot give precise estimates of *absolute* node ages



We need additional data (like fossils) to provide absolute time scale

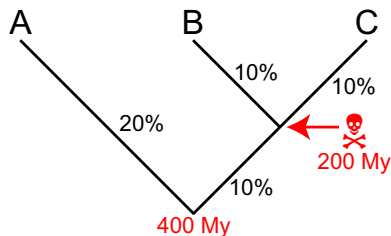


# CALIBRATING DIVERGENCE TIMES

Fossils (or other data) are necessary to estimate *absolute* node ages

There is **no information** in the sequence data for absolute time

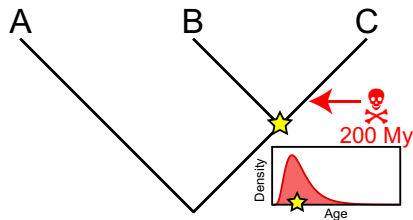
Uncertainty in the placement of fossils



# CALIBRATION DENSITIES

Bayesian inference is well suited to accommodating uncertainty in the age of the calibration node

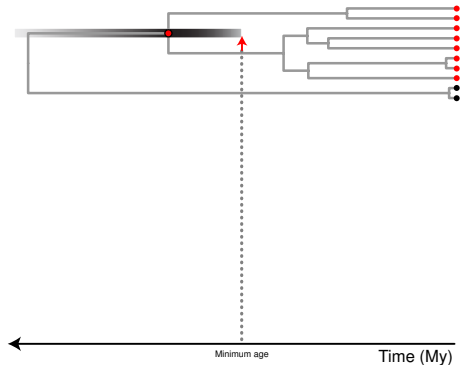
Divergence times are calibrated by placing parametric densities on internal nodes offset by age estimates from the fossil record



# FOSSIL CALIBRATION

Age estimates from fossils can provide **minimum** time constraints for internal nodes

Reliable **maximum** bounds are typically unavailable

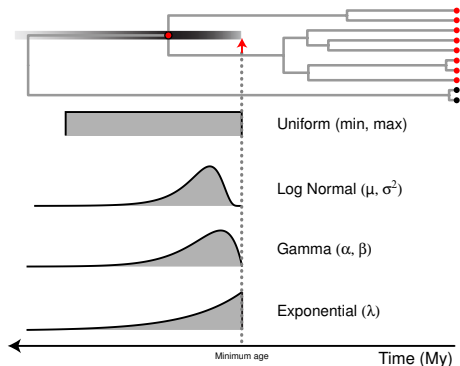


# DENSITIES ON CALIBRATED NODES

## Common practice in Bayesian divergence-time estimation:

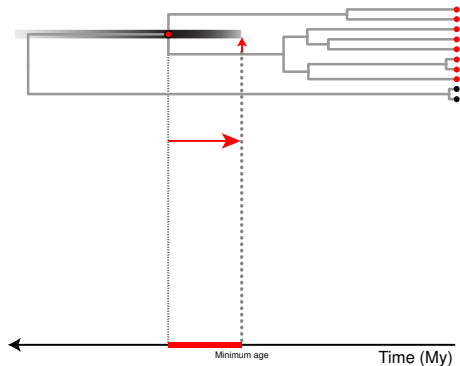
Parametric distributions are typically off-set by the age of the oldest fossil assigned to a clade

These prior densities do not (necessarily) require specification of maximum bounds



# DENSITIES ON CALIBRATED NODES

Typically interpreted as a “prior” on the waiting time between the divergence event and the age of the oldest fossil



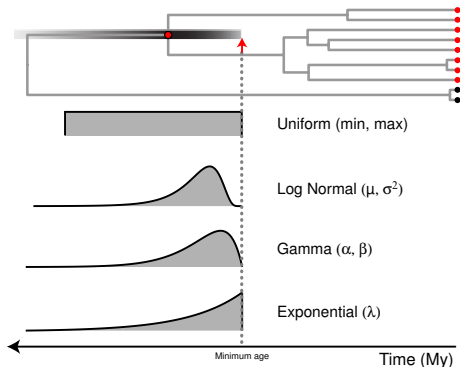


# DENSITIES ON CALIBRATED NODES

Common practice in Bayesian divergence-time estimation:

Estimates of absolute node ages are driven primarily by the calibration density

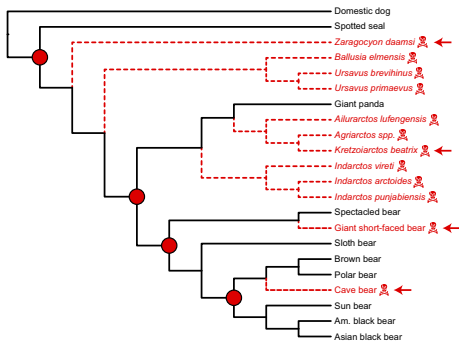
Specifying appropriate densities is a challenge for most molecular biologists



# IMPROVING FOSSIL CALIBRATION

We would prefer to eliminate the need for *ad hoc* calibration prior densities

Calibration densities do not account for diversification of fossils



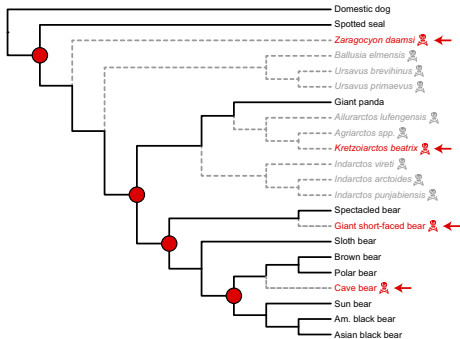
(Krause et al. *BMC Evol. Biol.* 2008; Abella et al. *PLoS ONE* 2012)

# IMPROVING FOSSIL CALIBRATION

We want to use all of the available fossils

## Example: Bears

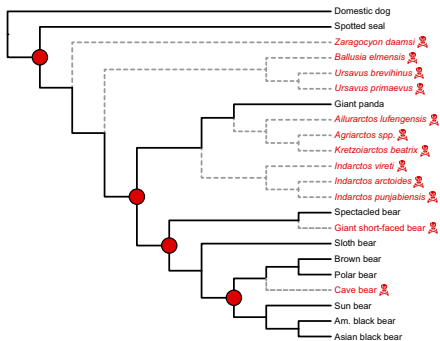
12 fossils are reduced to 4 calibration ages with calibration density methods



(Krause et al. *BMC Evol. Biol.* 2008; Abella et al. *PLoS ONE* 2012)

# IMPROVING FOSSIL CALIBRATION

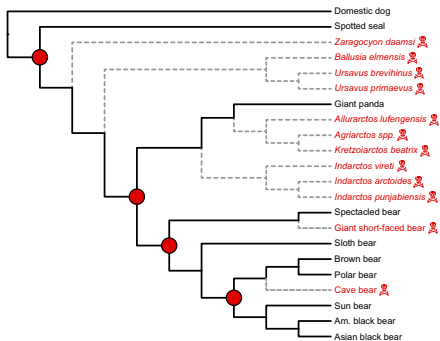
Because fossils are part of the diversification process, we can combine fossil calibration with birth-death models



(Krause et al. *BMC Evol. Biol.* 2008; Abella et al. *PLoS ONE* 2012)

# IMPROVING FOSSIL CALIBRATION

This relies on a branching model that accounts for **speciation, extinction, and rates of fossilization, preservation, and recovery**



(Krause et al. *BMC Evol. Biol.* 2008; Abella et al. *PLoS ONE* 2012)

# PALEONTOLOGY & NEONTOLOGY

“Except during the interlude of the [Modern] Synthesis, there has been limited communication historically among the disciplines of evolutionary biology, particularly between students of evolutionary history (paleontologists and systematists) and those of molecular, population, and organismal biology. **There has been increasing realization that barriers between these subfields must be overcome if a complete theory of evolution and systematics is to be forged.**”

Reaka-Kudla, M.L. & Colwell, R.: in E.C. Dudley (ed.), *The Unity of Evolutionary Biology: Proceedings of the Fourth International Congress of Systematic & Evolutionary Biology*, Discorides Press, Portland, OR, p. 16.



*Biology and Philosophy* **19**: 687–720, 2004.

© 2004 Kluwer Academic Publishers. Printed in the Netherlands.

## **The role of fossils in phylogeny reconstruction: Why is it so difficult to integrate paleobiological and neontological evolutionary biology?**

TODD GRANTHAM

*Department of Philosophy, College of Charleston, Charleston, SC 29424, USA*

*(e-mail: granthamt@cofc.edu)*

# COMBINING FOSSIL & EXTANT DATA

Statistical methods provide a way to integrate paleontological & neontological data

*Syst. Biol.* 50(6):913–925, 2001

## **A Likelihood Approach to Estimating Phylogeny from Discrete Morphological Character Data**

PAUL O. LEWIS

*Department of Ecology and Evolutionary Biology, The University of Connecticut, Storrs, Connecticut 06269-3043, USA;  
E-mail: paul.lewis@uconn.edu*

*Syst. Biol.* 61(6):973–999, 2012

© The Author(s) 2012. Published by Oxford University Press

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0>), which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

DOI:10.1093/sysbio/sys058

Advance Access publication on June 20, 2012

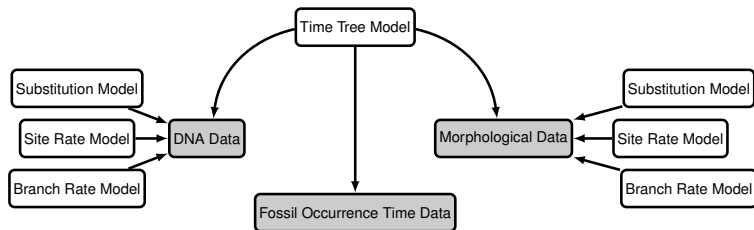
## **A Total-Evidence Approach to Dating with Fossils, Applied to the Early Radiation of the Hymenoptera**

FREDRIK RONQUIST<sup>1,\*</sup>, SERAINA KLOPFSTEIN<sup>1</sup>, LARS VILHELMSSEN<sup>2</sup>, SUSANNE SCHULMEISTER<sup>3</sup>, DEBRA L. MURRAY<sup>4</sup>, AND ALEXANDR P. RASNITSYN<sup>5</sup>



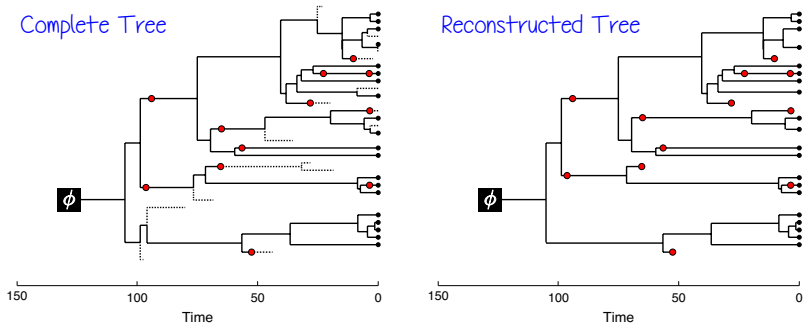
# COMBINING FOSSIL & EXTANT DATA

Combine models for sequence evolution, morphological change, & fossil recovery to jointly estimate the tree topology, divergence times, & lineage diversification rates



# MODELING THE TREE & OCCURRENCE TIMES

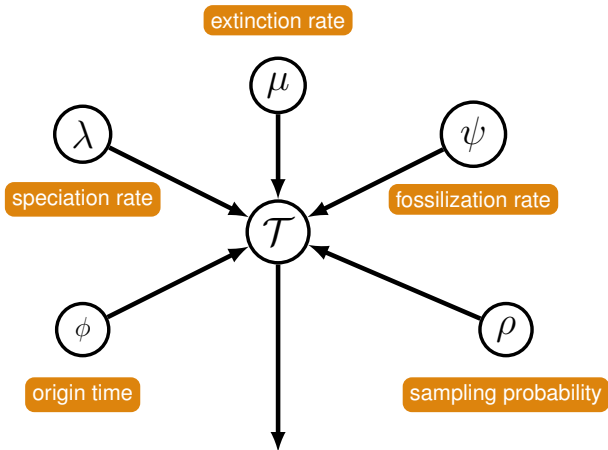
Stadler (2010) introduced a generating model for a serially sampled time tree — this is the *fossilized birth-death process*.



(Stadler. *Journal of Theoretical Biology* 2010)

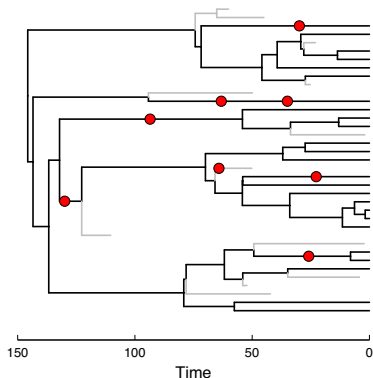
# PARAMETERS OF THE FBD

This graph shows the conditional dependence structure of the FBD model, which is a generating process for a sampled, dated time tree and fossil occurrences



# THE FOSSILIZED BIRTH-DEATH PROCESS (FBD)

Recovered fossil specimens provide historical observations of the diversification process that generated the tree of extant species



(Heath, Huelsenbeck, Stadler. *PNAS* 2014)

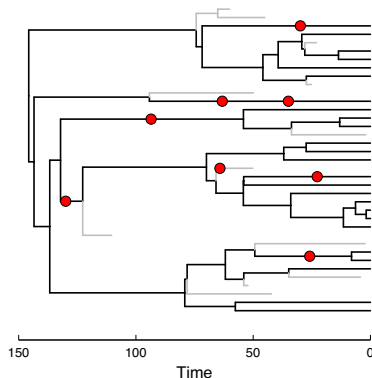
# THE FOSSILIZED BIRTH-DEATH PROCESS (FBD)

The probability of the tree and fossil observations under a birth-death model with rate parameters:

$\lambda$  = speciation

$\mu$  = extinction

$\psi$  = fossilization/recovery



(Heath, Huelsenbeck, Stadler. *PNAS* 2014)

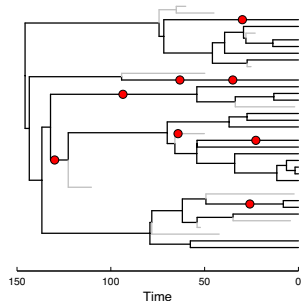
# SAMPLED ANCESTORS

Sampled lineages with sampled descendants

*Paleobiology*, 22(2), 1996, pp. 141–151

## On the probability of ancestors in the fossil record

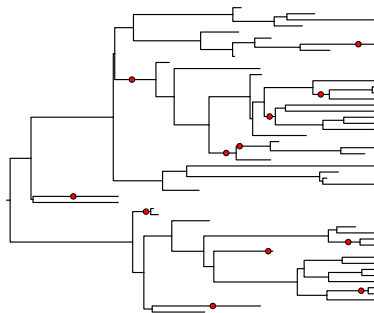
Mike Foote



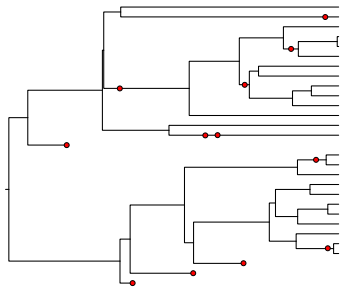
There is a non-zero probability of sampling ancestor-descendant relationships from the fossil record

# SAMPLED ANCESTORS

Complete FBD Tree



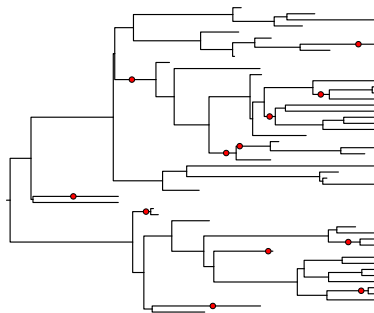
Reconstructed FBD Tree



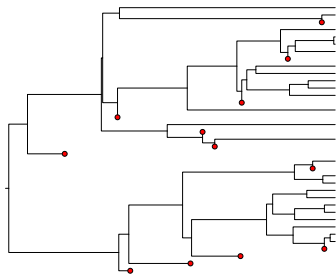
Because fossils & living taxa are assumed to come from a single diversification process, there is a non-zero probability of sampled ancestors

# SAMPLED ANCESTORS

Complete FBD Tree



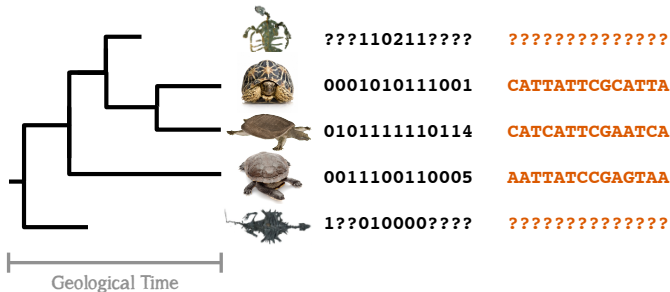
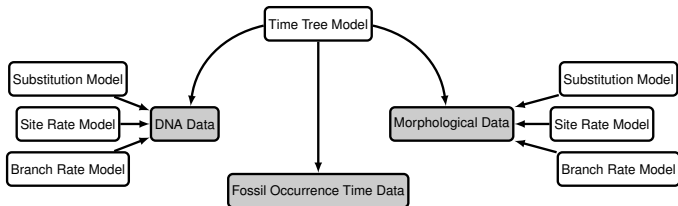
No Sampled Ancestor Tree



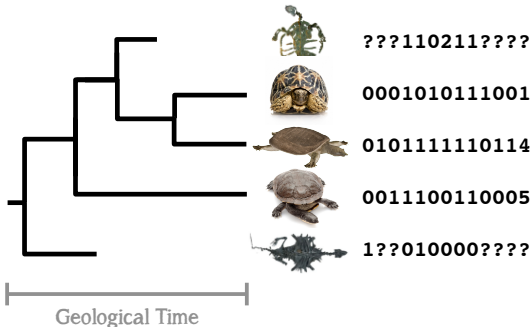
If all fossils are forced to be on separate lineages, this induces additional speciation events and will, in turn, influence rate & node-age estimates.



# COMBINING FOSSIL & EXTANT DATA



# MODELING MORPHOLOGICAL CHARACTER CHANGE



*Syst. Biol.* 50(6):913–925, 2001

## A Likelihood Approach to Estimating Phylogeny from Discrete Morphological Character Data

PAUL O. LEWIS

*Department of Ecology and Evolutionary Biology, The University of Connecticut, Storrs, Connecticut 06269-3043, USA;  
E-mail: paul.lewis@uconn.edu*

# MODELING MORPHOLOGICAL CHARACTER CHANGE

## The Lewis Mk model

Assumes a character can take  $k$  states

Transition rates between states are equal (symmetric)

$$Q = \alpha \begin{bmatrix} 1 - k & 1 & \dots & 1 \\ \vdots & 1 - k & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 - k \end{bmatrix}$$

**T1**    **0**

**T2**    **0**

**T3**    **1**

**T4**    **2**

**T5**    **2**

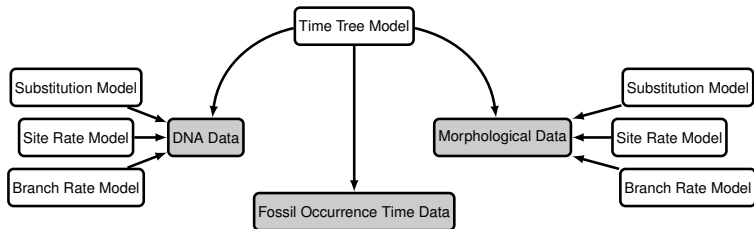
**T6**    **1**

**T7**    **1**

(Lewis. *Systematic Biology* 2001)

# "TOTAL-EVIDENCE" ANALYSIS

Integrating models of molecular and morphological evolution with improved tree priors enables joint inference of the tree topology (extant & extinct) and divergence times

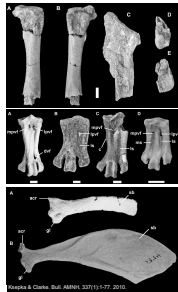


# PENGUIN DIVERSITY IN DEEP TIME

How does our understanding of penguin evolution improve when we consider both extant and fossil taxa?

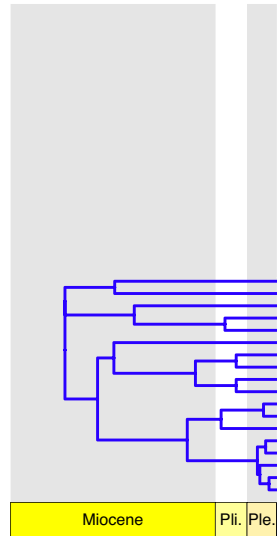


"Penguin Party" by Kate Dzikiewicz

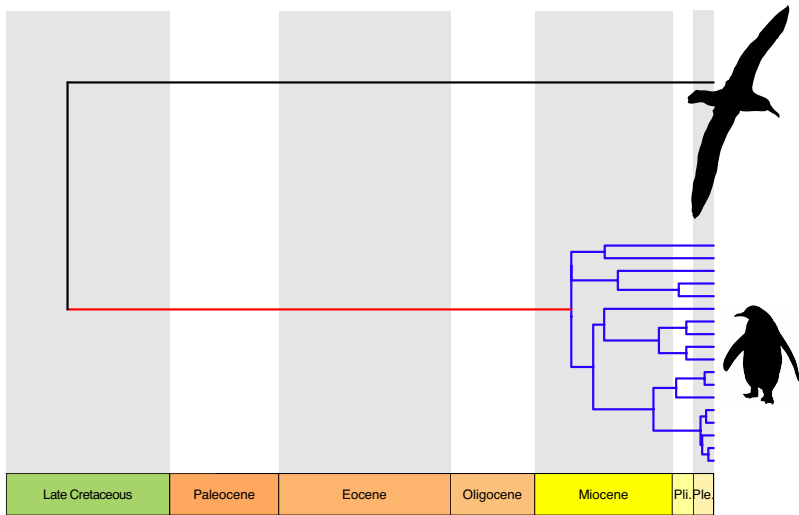


Artistic reconstructions by: Stephanie Abramowicz for Scientific American  
Fonyo, R.E. and D.T. Ksepka. The Strangest Bird Scientific American 307, 56 – 61 (2012)

# PENGUIN DIVERSITY



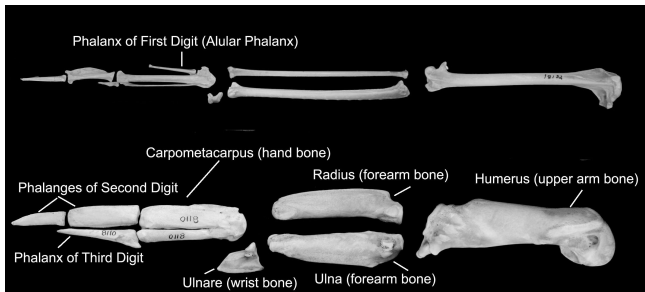
# PENGUIN DIVERSITY



(silhouette images from <http://phylopic.org>)

# WHAT MAKES A PENGUIN A PENGUIN?

Flattened, solid wing-bones



(image courtesy of D. Ksepka <https://fossilpenguins.wordpress.com>)

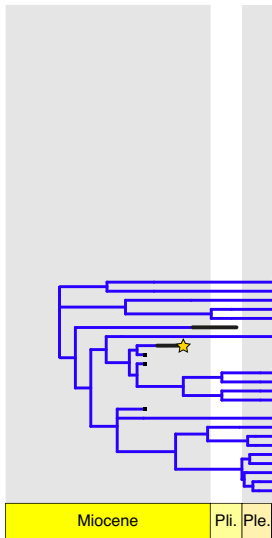


# FOSSIL PENGUIN DIVERSITY

*Spheniscus urbinai*  
Holotipo MUSM 401



Martin Chávez



(*S. urbinai* holotype fossil, 5-7 MYA, image by Martin Chávez)

# PENGUINS IN THE OLIGOCENE

## *Kairuku*

- ~1.5 m tall
- slender, with narrow bill
- scapula & pygostyle are more similar to non-penguins
- ~27 Mya

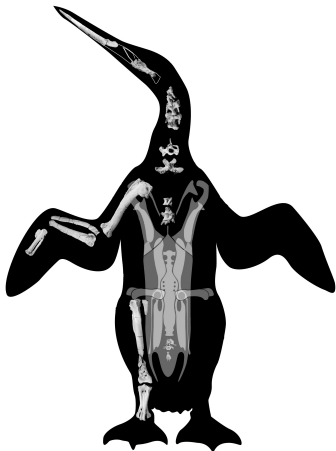


(Ksepka, Fordyce, Ando, & Jones, *J. Vert. Paleo.* 2012)

# PENGUINS IN THE PALEOCENE

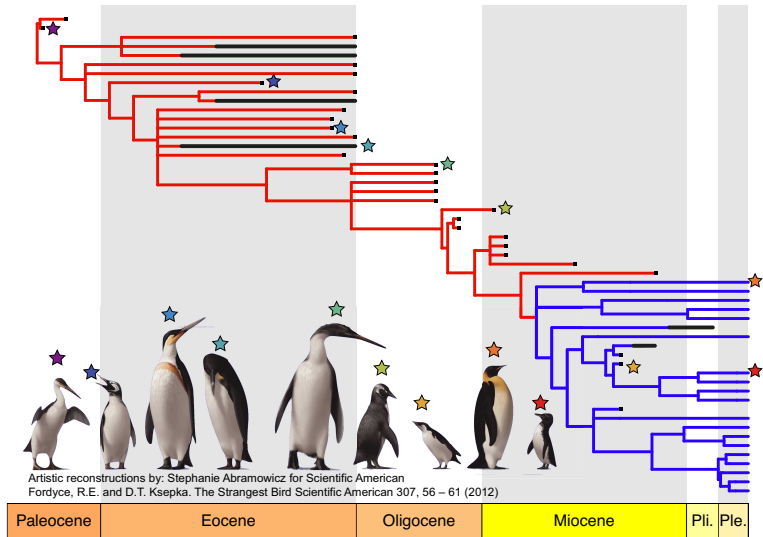
## *Waimanu*

- oldest known penguin species
- intermediate wing morphology
- ~58–61.6 Mya

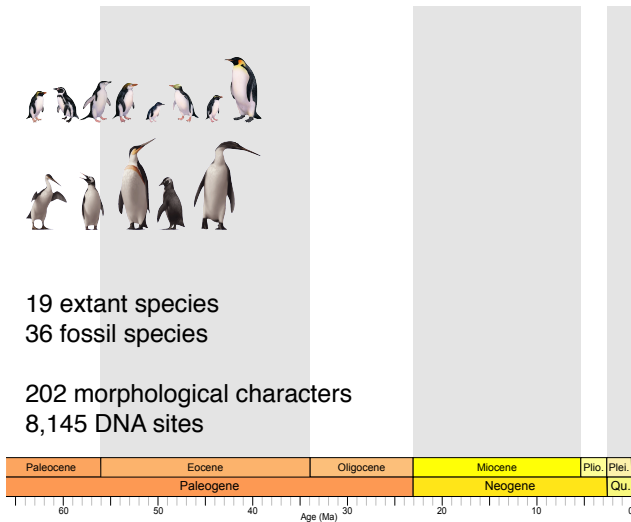


(Slack et al., *Mol. Biol. Evol.* 2006)

# PENGUIN DIVERSITY IN DEEP TIME



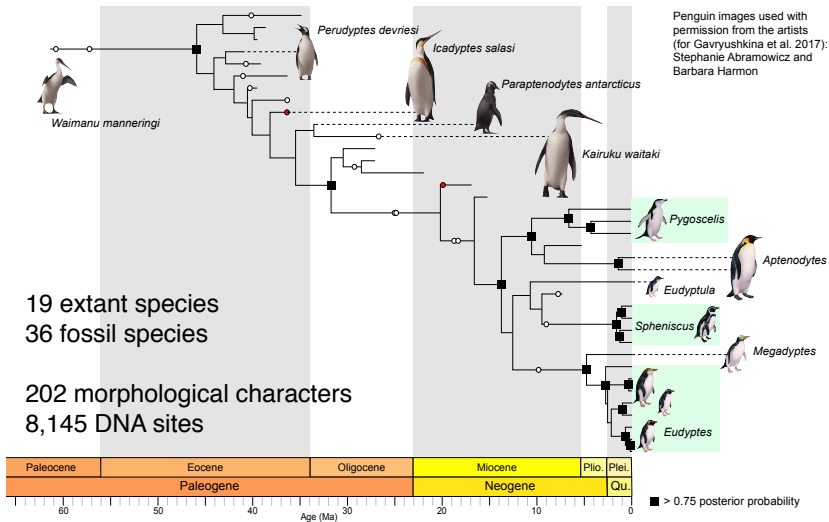
# PENGUIN DIVERSITY IN DEEP TIME



Penguin images used with permission from the artists (for Gavryushkina et al. 2017): Stephanie Abramowicz and Barbara Harmon

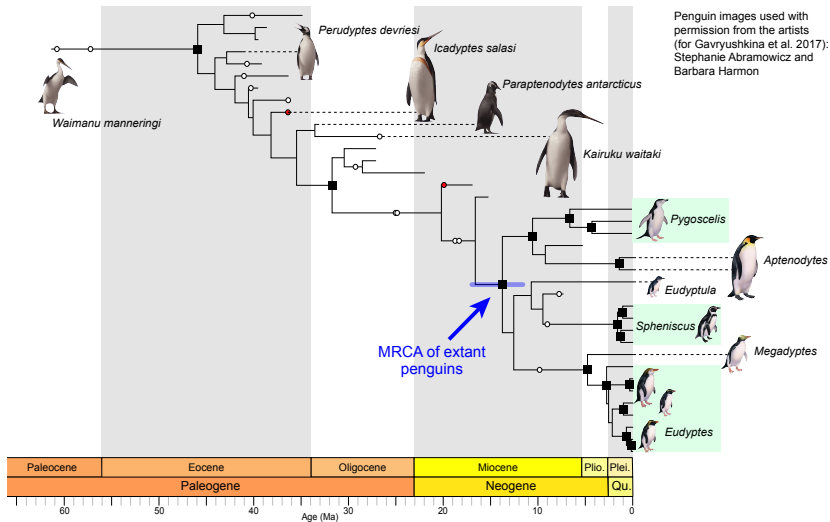
(Gavryushkina, Heath, Ksepka, Welch, Stadler, Drummond. 2017. *Syst. Biol.*)

# PENGUIN DIVERSITY IN DEEP TIME



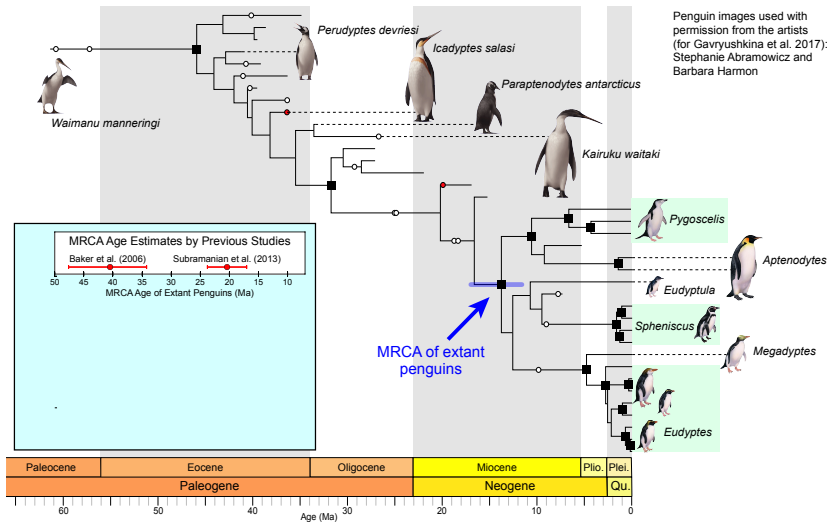
(Gavryushkina, Heath, Ksepka, Welch, Stadler, Drummond. 2017. *Syst. Biol.*)

# PENGUIN DIVERSITY IN DEEP TIME



(Gavryushkina, Heath, Ksepka, Welch, Stadler, Drummond. 2017. *Syst. Biol.*)

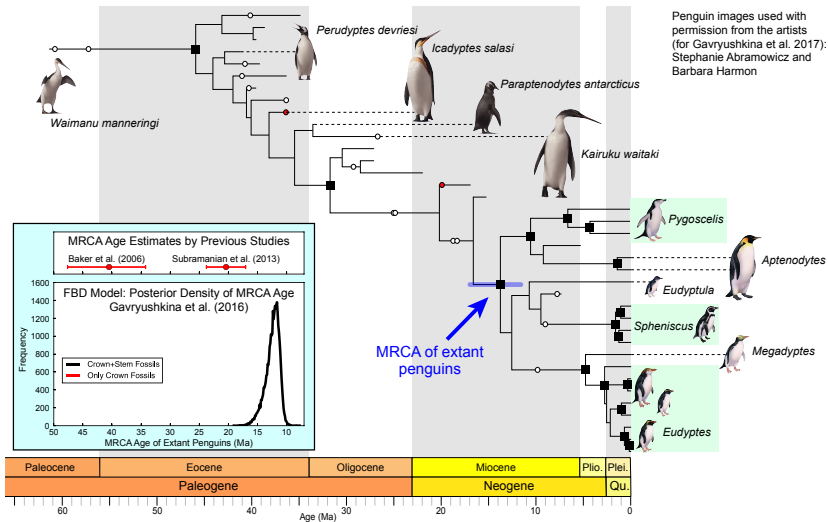
# PENGUIN DIVERSITY IN DEEP TIME



(Gavryushkina, Heath, Ksepka, Welch, Stadler, Drummond. 2017. *Syst. Biol.*)



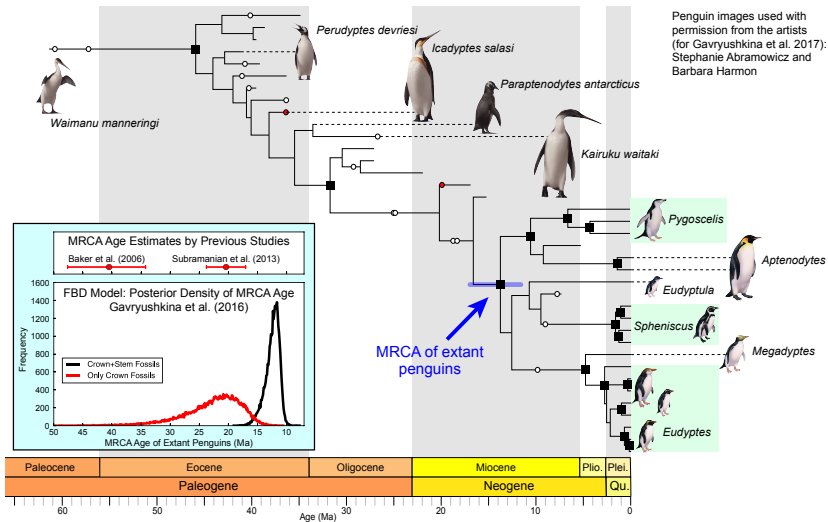
# PENGUIN DIVERSITY IN DEEP TIME



Penguin images used with permission from the artists (for Gavryushkina et al. 2017): Stephanie Abramowicz and Barbara Harmon

(Gavryushkina, Heath, Ksepka, Welch, Stadler, Drummond. 2017. *Syst. Biol.*)

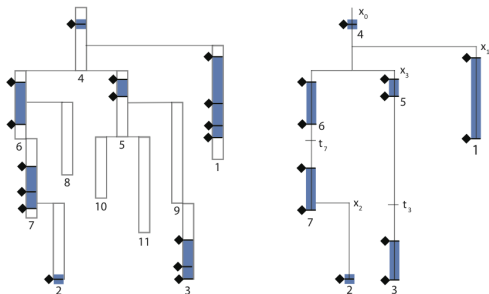
# PENGUIN DIVERSITY IN DEEP TIME



(Gavryushkina, Heath, Ksepka, Welch, Stadler, Drummond. 2017. *Syst. Biol.*)

# Fossil DATA & PHYLOGENIES

Through collaboration with paleontologists, we are building models to account for the structure of the fossil record and the nature of paleontological data



The fossilized birth-death model for the analysis of stratigraphic range data under different speciation modes

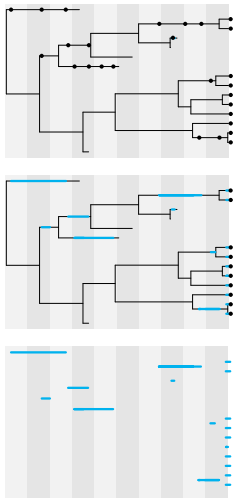
Tanja Stadler<sup>a,b,\*</sup>, Alexandra Gavryushkina<sup>a,b</sup>, Rachel C.M. Warnock<sup>a,b</sup>,  
Alexei J. Drummond<sup>c</sup>, Tracy A. Heath<sup>d</sup>

*Journal of Theoretical Biology* 447:41-55 (2018)

# FOSSIL DATA & PHYLOGENIES

The FBD model can accommodate different kinds of paleontological data

- specimen-level sampling
- when the fossil data are only coded for first and last occurrences (stratigraphic ranges)
- when only stratigraphic range data are available



(figure courtesy of R. Warnock)

# FBD FOR STRATIGRAPHIC RANGE DATA

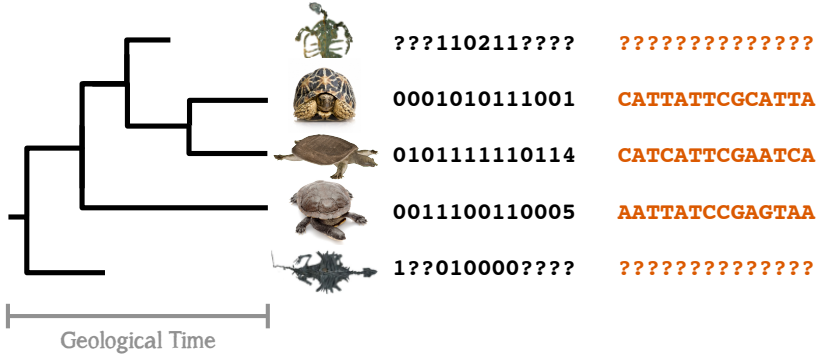
Estimate rates of speciation, extinction, & fossil recovery when no phylogenetic data are available



(figure from <https://gerardofurtado.com/sr/sr.html>, using data from Sepkoski 2002)

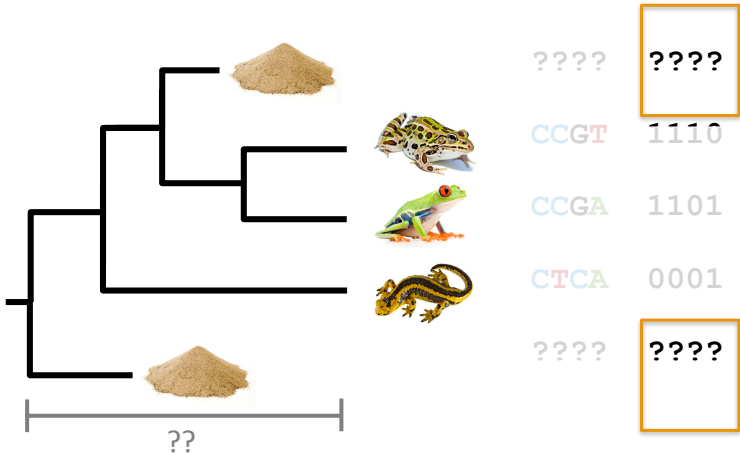
The FBD model explicitly accounts for incomplete species sampling, as well as uncertainty in speciation and extinction times and the phylogeny

# MOLECULES + MORPHOLOGY + FOSSILS

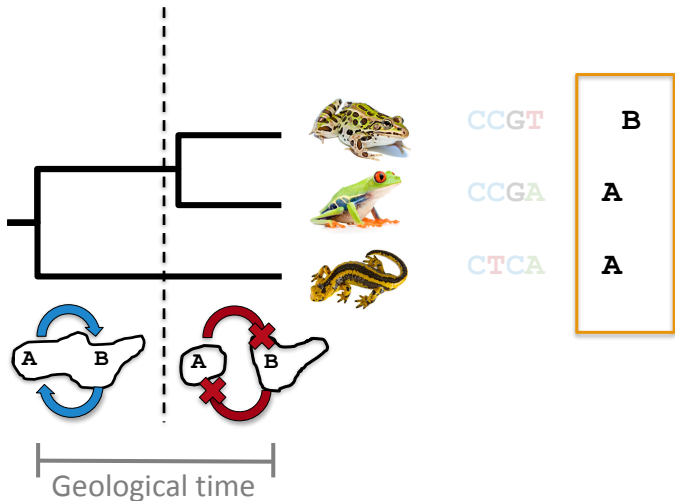


(based on slides by M. Landis)

# ...but I study amphibians...



# Molecules + biogeography + paleogeography



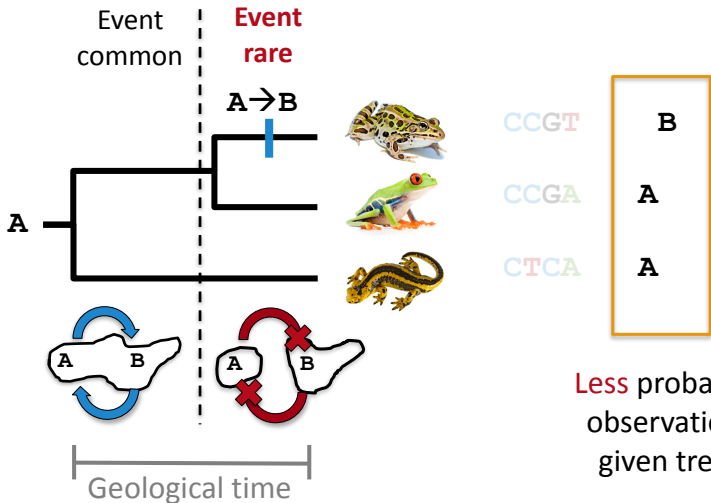
**+ Paleogeography**

Landis, 2016

(slides courtesy of M. Landis, <http://bit.ly/2aHqB4>)



# Events should occur *before* areas split

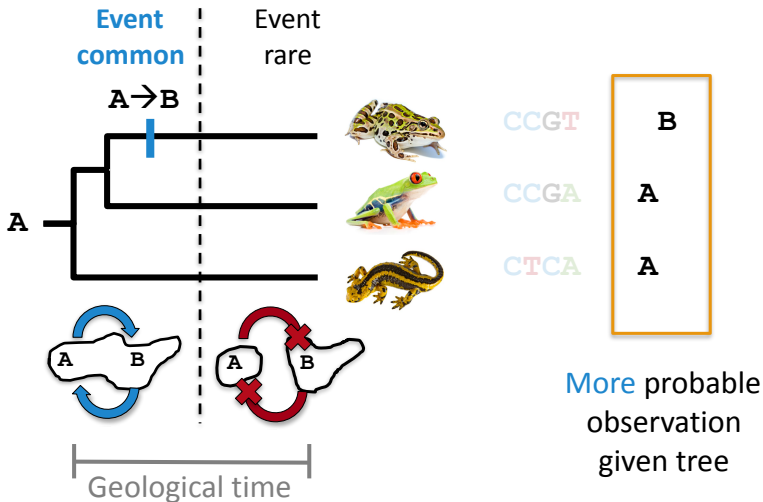


+ Paleogeography

Landis, 2016

(slides courtesy of M. Landis, <http://bit.ly/2aHqB4>)

# Events should occur *before areas split*

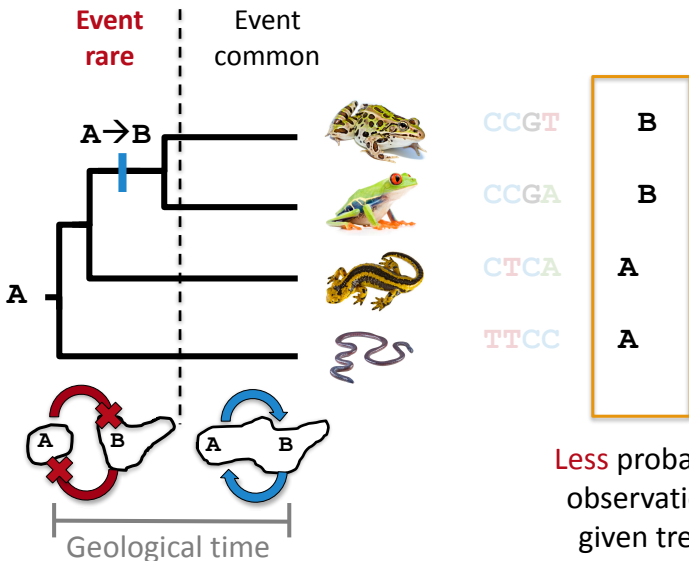


+ Paleogeography

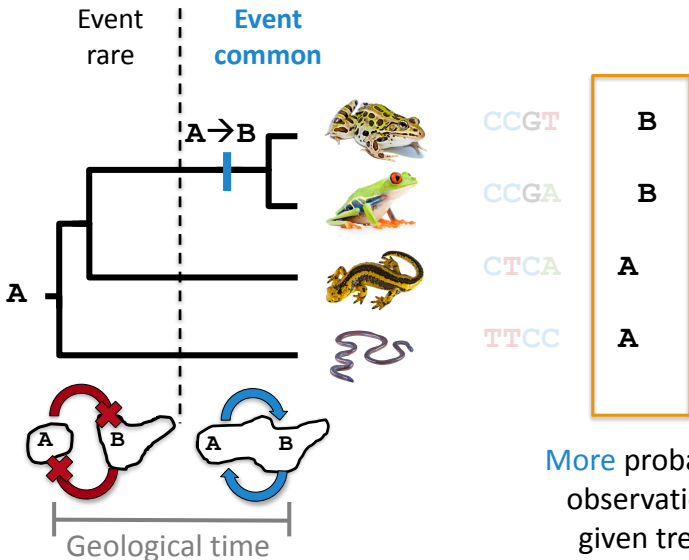
Landis, 2016

(slides courtesy of M. Landis, <http://bit.ly/2aHqB4>)

# Events should occur *after* areas merge



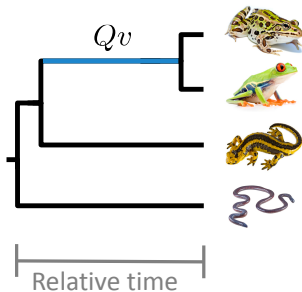
# Events should occur *after* areas merge



# BIOGEOGRAPHIC DATING

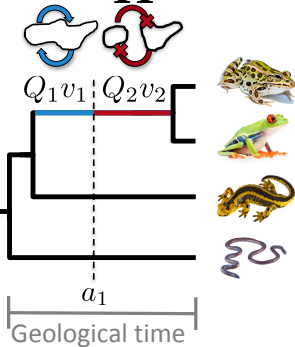
## Constant Model

$$P_{ij}(v) = \exp\{Qv\}$$



## Epoch Model

$$P_{ij}(\mathbf{v}) = \prod \exp\{Q_k v_k\}$$

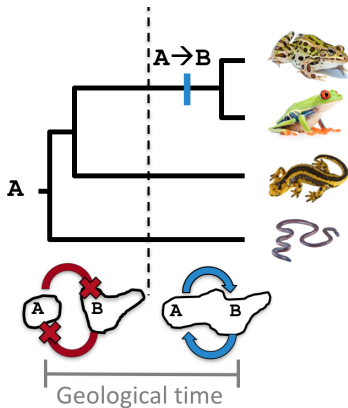


**Epoch model** Ree et al., 2005  
Bielejec et al., 2014

# BIOGEOGRAPHIC DATING

## Fossil-free calibration

- data: molecular sequences
- data: biogeographic ranges
- empirical paleogeographic model that alters the rates of biogeographic change over time

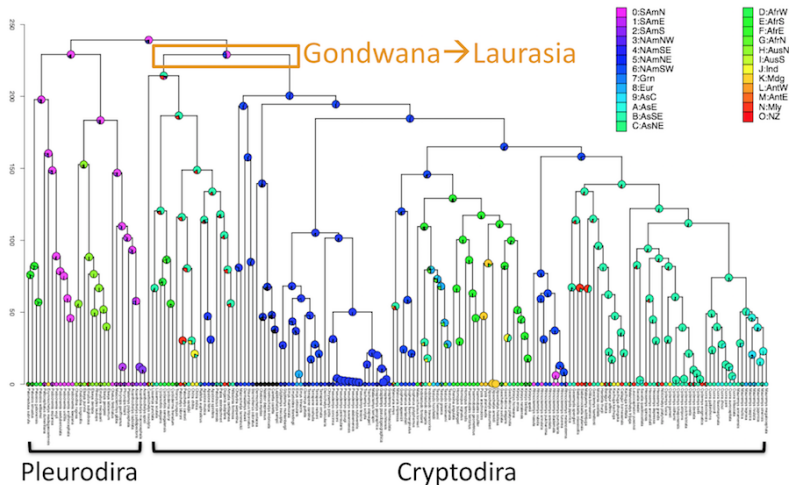


Landis. 2017. "Biogeographic Dating of Speciation Times Using Paleogeographically Informed Processes". *Systematic Biology*

doi: [10.1093/sysbio/syw040](https://doi.org/10.1093/sysbio/syw040).

# DATING + ANCESTRAL AREA RECONSTRUCTION

## Ancestral area estimates (+G)



# BAYESIAN ANALYSIS UNDER COMPLEX MODELS

The models I have described so far are primarily available in two software packages

## ***RevBayes***

[revbayes.com](http://revbayes.com)

flexible model construction using graphical models and an interpreted language for modeling, simulation, and Bayesian inference in evolutionary biology, particularly phylogenetics

tutorials:

[revbayes.com/tutorials](http://revbayes.com/tutorials)

## ***BEAST 2***

[beast2.org](http://beast2.org)

a modular program for using Bayesian inference to infer rooted trees for addressing questions in population genetics, phylogenetics, epidemiology, and macroevolution

tutorials:

[taming-the-beast.org](http://taming-the-beast.org)



# LARGE DATASETS & BAYESIAN INFERENCE

Alignments with thousands of loci and/or thousands of taxa are not practical with fully Bayesian methods, particularly under complex, hierarchical models.

Large datasets lead to long mixing times because the number of parameters requires many proposals in order to effectively sample the joint posterior density

## *What to do?*

- Simplify the problem (in ways that do not also sacrifice accuracy)



Pecos

# APPROXIMATE COMPUTATION OF THE LIKELIHOOD

Computing the full likelihood of an alignment on a tree is computationally intensive

Bayesian inference under complex models from large datasets requires the likelihood to be calculated millions of times

The practical solution is to compute the likelihood faster, which is possible if we can get a reasonable approximation

## **Approximate Likelihood Calculation on a Phylogeny for Bayesian Estimation of Divergence Times**

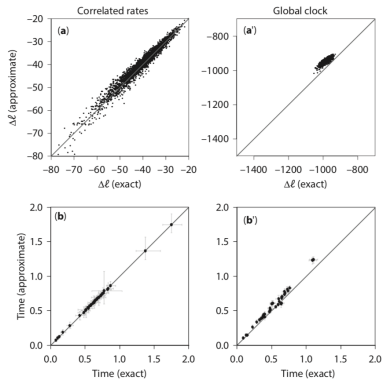
Mario dos Reis, Ziheng Yang ✉

*Molecular Biology and Evolution*, Volume 28, Issue 7, 1 July 2011, Pages 2161–2172,

# APPROXIMATE COMPUTATION OF THE LIKELIHOOD

A 2-step analysis:

- 1 branch lengths are first estimated on the unrooted tree along with elements needed for approximation method
- 2 estimate divergence times using MCMC, with inputs from 1 for calculating the approximate likelihood for each M-H move



**Fig. 10.14** (a and a') The log likelihood (or the difference from the highest log likelihood at the MLEs) and (b and b') posterior estimates (means and 95% CIs) of divergence times calculated using the exact and approximate likelihood methods under the relaxed clock and strict clock models. The mitochondrial data for 36 mammalian species are analysed. Redrawn from the plots for the arcsine transform in Figures 3 and 5 of dos Reis and Yang (2011).

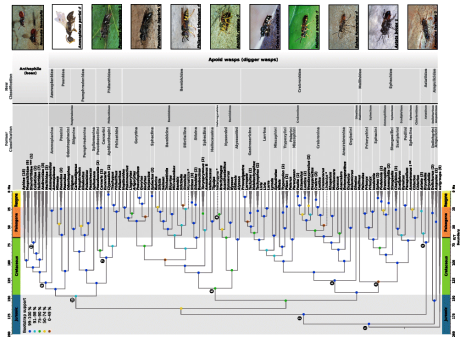
(figure from Yang, 2014. *Molecular Evolution: A Statistical Approach*)

# BAYESIAN DATING WITH BIG(GER) DATA

Using approximate likelihood computation in MCMCTree

## *Apoïd Wasps & Bees*

- 174 taxa
- 284,607 bp (195 single-copy protein-coding genes)
- 4 fossil calibrations



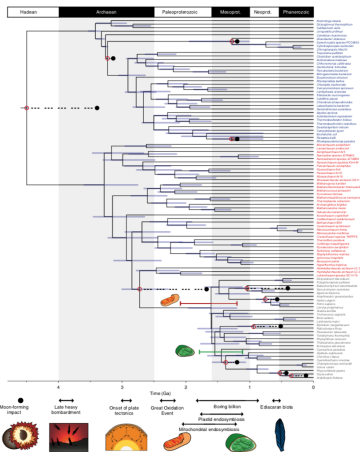
Phylogenomic analysis of Apoidea sheds new light on the sister group of bees (Sann et al. 2018. *BMC Evol. Biol.*)

# BAYESIAN DATING WITH BIG(GER) DATA

Using approximate likelihood computation in MCMCTree

## Life

- 102 taxa
- 14,745 bp (29 universally distributed, protein-coding genes)
- 11 fossil calibrations



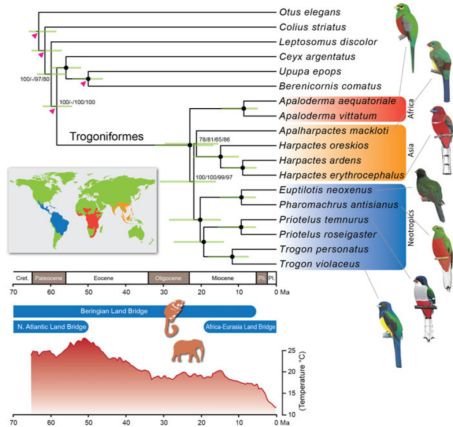
Integrated genomic and fossil evidence illuminates life's early evolution and eukaryote origin (Betts et al. 2018. *NE&E*)

# BAYESIAN DATING WITH BIG(GER) DATA

Using approximate likelihood computation in MCMCTree

## Trogon

- 12 ingroup + 6 outgroup taxa
- 2.9 Million bp (3581 UCE loci)
- 4 fossil calibrations



Rapid Laurasian diversification of a pantropical bird family during the Oligocene-Miocene transition (Oliveros et al. in press. *Ibis*)